

ChatGPT

Internals for coders

—— Jakub Švehla ——

Goals

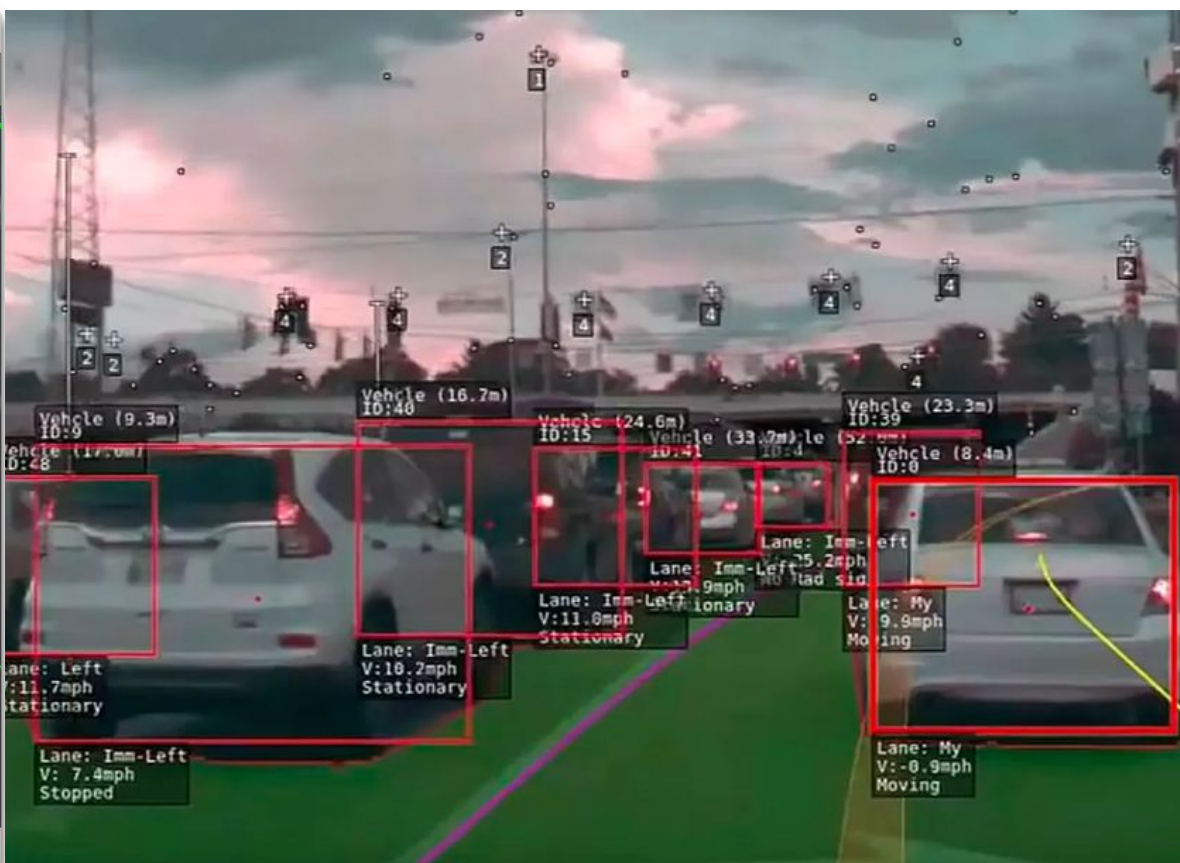
1. What is generative AI?
2. Understand how chatGPT is implemented
3. Current challenges in the world of Natural Language Processing
4. Ezoteric future chatGPT predictions

What is generative AI?

AI

- Image Classifiers
- User Clustering & auto segmentation in Online Advertising
- Self-organisation & self-assemblings
- **Generative models**
 - **Images**
 - **3D models**
 - **Texts (NLP - natural language processing)**
- ...

Classifiers



Self assembling AI

https://sebastianrisi.com/self_assembling_ai/

Generative AI - Images



Generative AI - 3D objects



Generative AI - Natural language processing

```
// write a binary search algorithm
const binarySearch = (arr, target) => {
  let left = 0;
  let right = arr.length - 1;
  let middle = Math.floor((left + right) / 2);
  while (arr[middle] !== target && left <= right) {
    if (target < arr[middle]) {
      right = middle - 1;
    } else {
      left = middle + 1;
    }
    middle = Math.floor((left + right) / 2);
  }
  return arr[middle] === target ? middle : -1;
}
```



write a story about myself working in Bloomreach



Once upon a time, there was a talented young professional named Alex who dreamed of working for a leading tech company that would challenge and inspire him every day. After applying to numerous positions and undergoing a rigorous interview process, Alex was thrilled to receive an offer from Bloomreach, one of the most innovative and dynamic companies in the tech industry.

From day one, Alex was impressed by the culture of Bloomreach, which fostered collaboration, creativity, and continuous learning. As a member of the product development team, Alex was tasked with designing and implementing new features for Bloomreach's groundbreaking search and merchandising platform. Although the work was challenging,

How computers see images



62	62	63	64	65	66	67	67	69	70	71	72	72	73	73	73	72	72	71	70	69	67	66	66	65	63	62	61	60	6		
61	62	63	64	66	66	67	68	68	69	70	71	71	72	72	73	72	72	71	71	70	69	68	66	66	65	63	62	61	60	6	
61	62	63	64	66	66	68	68	69	70	70	71	72	73	73	73	72	71	71	69	68	67	66	66	65	65	64	63	62	61	6	
61	63	64	64	66	67	68	68	68	69	70	71	71	73	73	74	73	73	73	71	70	69	68	66	66	65	64	63	62	61	61	6
61	63	64	65	67	68	69	69	70	70	71	71	72	55	53	69	72	72	71	71	70	69	68	67	66	65	64	63	62	60	60	6
63	64	65	66	67	68	69	69	70	70	71	72	42	4	5	11	58	72	71	71	69	69	68	67	66	65	64	62	62	60	59	5
63	65	66	66	68	68	69	70	71	71	71	72	18	4	4	7	8	66	71	70	69	68	68	67	66	65	64	63	61	59	59	5
63	65	67	67	68	69	69	70	71	71	72	64	4	27	24	54	33	20	52	64	68	68	67	66	65	64	63	62	61	59	58	5
64	65	66	66	68	69	70	71	41	24	24	12	17	24	48	60	37	43	30	52	66	68	67	66	65	64	63	61	60	59	58	5
65	66	67	67	68	69	71	40	6	6	6	5	34	36	12	47	34	17	20	54	43	63	67	66	65	64	63	62	60	59	58	5
64	65	66	66	68	69	38	6	6	5	5	7	16	19	4	47	44	27	24	40	67	66	66	65	65	64	63	61	60	59	58	5
63	64	65	65	67	30	6	6	5	5	5	6	8	9	20	27	51	78	41	44	66	65	65	65	65	64	63	62	60	59	58	5
63	64	65	65	34	5	5	5	5	5	5	5	4	19	6	7	54	64	20	59	65	65	64	64	64	63	62	61	60	59	57	5
63	64	64	65	14	5	6	5	5	4	5	4	18	7	5	4	19	10	11	65	64	64	64	63	61	66	62	61	60	59	58	5
63	64	64	65	53	7	4	5	6	6	7	10	6	5	5	4	21	24	18	64	64	64	63	62	64	65	62	62	60	59	58	5
64	64	64	64	65	50	4	4	4	5	11	16	6	6	4	6	35	16	20	66	64	64	63	61	72	67	63	62	61	59	58	5
64	64	64	64	65	46	4	4	4	5	6	9	8	5	29	10	43	56	20	57	64	64	63	61	70	67	62	64	65	59	59	5
64	64	64	65	66	27	5	4	4	5	6	6	6	18	66	20	57	60	46	36	75	70	62	61	70	67	62	61	60	59	58	5
49	50	62	65	57	5	5	6	5	6	6	6	6	51	59	23	60	58	44	22	63	71	72	60	69	68	61	60	58	59	59	5
42	52	57	52	26	5	5	5	5	5	5	5	5	70	50	43	61	62	64	39	42	64	60	62	56	63	65	65	67	61	53	5
32	32	32	33	6	5	5	5	5	5	6	6	11	39	21	33	51	50	45	46	18	32	36	33	23	44	70	71	51	42	27	3
50	50	51	39	5	5	5	5	6	5	6	6	42	69	28	34	42	39	43	37	26	29	40	26	29	26	35	42	35	33	18	1
52	53	51	22	5	5	5	5	6	5	6	5	44	56	17	51	54	53	54	56	51	22	54	54	55	55	54	53	53	53	52	5
54	54	53	8	5	5	5	5	6	5	6	13	52	42	21	51	54	51	49	49	50	22	41	45	42	42	41	40	41	44	43	4
52	52	54	36	8	5	5	6	6	5	6	28	55	32	32	54	53	51	51	51	51	44	25	51	51	49	49	50	49	48	46	4
54	54	52	53	30	7	5	6	6	5	6	40	54	29	52	51	53	56	55	52	52	51	38	52	52	50	49	46	46	45	46	4
51	52	51	53	27	14	5	4	5	4	7	47	51	21	39	49	47	49	52	52	49	35	31	48	46	47	47	47	47	46	46	4
48	50	51	53	25	14	17	8	4	4	17	46	40	18	43	47	46	49	52	54	53	53	54	18	50	49	46	47	47	47	47	4
49	49	49	49	22	12	20	24	6	14	35	51	39	48	48	50	51	51	49	51	51	52	50	41	58	48	47	47	47	45	45	4
51	49	50	50	22	13	19	36	13	12	42	50	40	73	50	50	50	49	48	49	48	49	45	51	46	44	44	44	42	45	4	
47	49	49	47	20	16	26	39	21	15	36	48	42	61	47	48	51	47	50	51	51	51	49	47	47	52	47	47	44	43	45	4

How computers see text

Computers are able to `/` `*` `+` `-`

Can't do math between those number

Hello Bloomreach!	72 101 108 108 111 32 66 108 111 111 109 114 101 97 99 104 33 10
where is the closest pizza place	119 104 101 114 101 32 105 115 32 116 104 101 32 99 108 111 115 101 115 116 32 112 105 122 122 97 32 112 108 97 99 101

Reimplement ChatGPT from scratch



LLM (little language model)

(Chat Little Generative pretrained transformerish)

1. Download Corpus
2. Tokenizer
3. One hot encoding
4. Word Embeddings
5. Contextual Word Embeddings
6. Next word prediction
7. Fine tuning for Instruct GPT & ChatGPT & AI safeness
8. Scaling hardware

LLM Corpus

Bloomreach Loves Bloomreach

Loves!!!!

Bloomreach!!!

Loves Bloomreach

Bloomreach Loves Bloomreach !!!

! Bloomreach Loves Bloomreach !

Tokenizer

Tokenizer

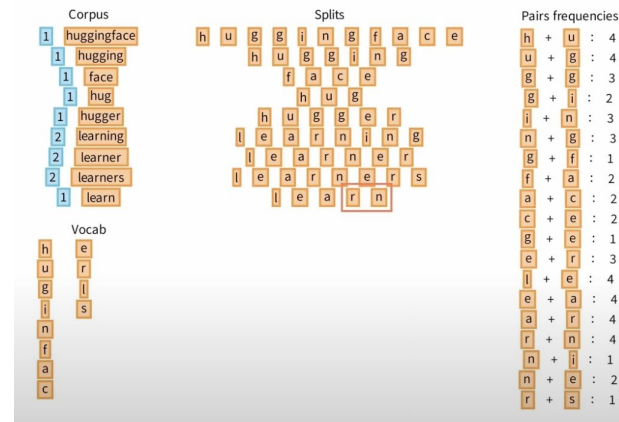
- Split corpus into smaller tokens (f.e: words)
- How many possible tokens should we support?
- Compute tokens from our text corpus

Tokenizer problematic exceptions

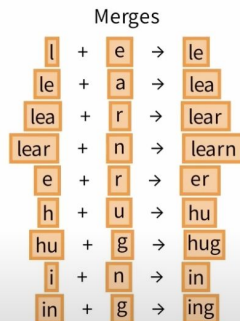
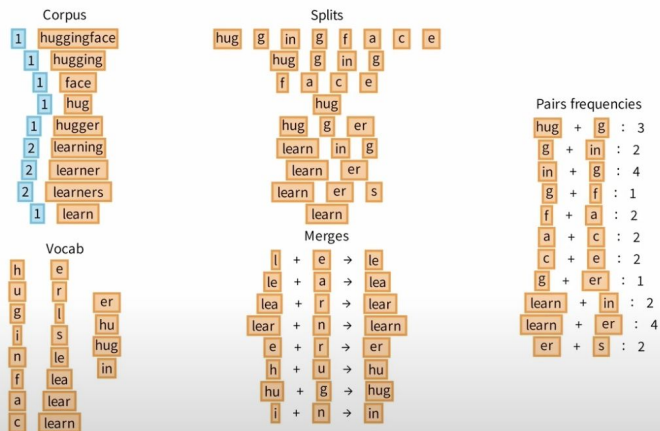
Typos	bloormeach	
Case sensitive	BLOOMREACH	→ [BLOOMREACH] [bloomreach]
Numbers	1234	→ [1 2 3 4] [1234]
Abbreviation	U.S.A.	→ [U . S . A .] [U.S.A.]
idiomatic expressions	PCI Express	→ ["PCI Express"] ["PCI" "Express"]
New Words	GPT	→ [GPT] [G P T]
说/説	说/説	→ [说 / 説]

Tokenizer - BPE

- Compute probabilities (relevancy) of sequences
- <https://platform.openai.com/tokenizer> (tik-token)



this is the hugging face course . this chapter is about
tokenization . this section shows several tokenizer
algorithms .



3x this 2x is 1x the
1x hugging 1x face
1x course 3x . 1x chapter
1x about 1x tokenization
1x section 1x shows
1x several 1x tokenizer
1x algorithms

Bloomreach	Loves	Bloomreach			
Loves	!	!	!	!	
Bloomreach	!	!	!		
Loves	Bloomreach				
Bloomreach	Loves	Bloomreach	!	!	!
!	Bloomreach	Loves	Bloomreach	!	

One Hot Encoding

One Hot Encoding

- Transform a set of “non continuous” data into 2D points
- Create a new dimension for each token
- Easy for human readable and setting the data

(Bloomreach, Loves, !)						
[2, 1, 0]	Bloomreach	Loves	Bloomreach			
[1, 1, 4]	Loves	!	!	!	!	
[1, 0, 3]	Bloomreach	!	!	!		
[1, 1]	Loves	Bloomreach				
[2, 1, 3]	Bloomreach	Loves	Bloomreach	!	!	!
[2, 1, 2]	!	Bloomreach	Loves	Bloomreach	!	

```

lovely(sentence1) > lovely(sentence3)
lovely([2, 1, 0]) > lovely([1, 0, 3])
1 > 0 = true

```

One Hot Encoding

Props

- Simple to create
- Possible to do arithmetics
- Simple to understand
- Used by many algorithms like
 - BM25 / TF-IDF (term frequency–inverse document frequency)

Cons

- Can't recognize synonyms opposite words
- In the vectors we do not have meaning
- Gigantic count of dimensions (tens of thousands)
- The more tokens we have the more dimensions we need to represent basic words.

INTERRUPTION!!!

Math basics

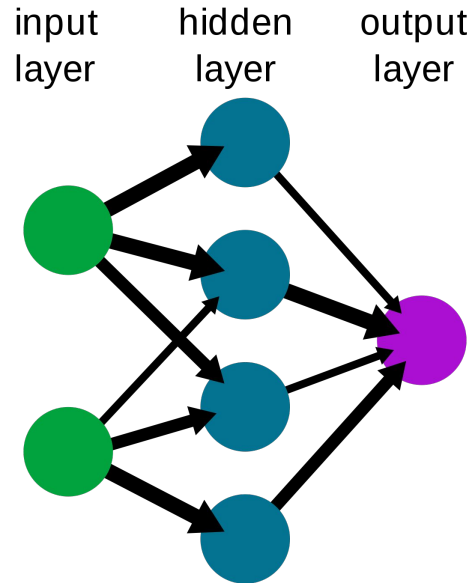
Math int - Math operations

1. Vector
2. Vector space
3. Linear transformation
4. Non-linear transformation
5. Neural networks

Math int - Neural networks

A program that helps to find the best matrices that we will apply to our data to approximate general-purpose behavior.

A simple neural network



Word embeddings

Word Embeddings

- Dimension-reduced One hot encoded vector space (latent space)
- Vectorized tokens with meanings in the latent space
- The latent space may be multilingual
- Each dimension describes some feature of the data

Word Embeddings - training

- Word2Vec + GloVE (most famous algorithms)
- Unsupervised learning on the giants corpuses
- What does word mean
- Looking for the most universal linear transformation
 - lower embedding dimensions
 - Middle word masking

Word Embeddings

Props

- Lookup table $O(1)$ runtime complexity
- Can do arithmetics
- Measure similarity
 - dot product
 - cosine similarity
 - euclidean distance)
- MLM (multi languages models) Transfer learning
 - Abstraction over natural languages ($\text{vec}(\text{car}) \sim \text{vec}(\text{auto})$)
 - F.e.: Need training data in english to learn some stuffs in czech

Cons

- No context word by meaning
- No sentence embeddings
- No positional encoding

TOKENIZER

ONE HOT
ENCODING

WORD
embeddings
 $O(1)$

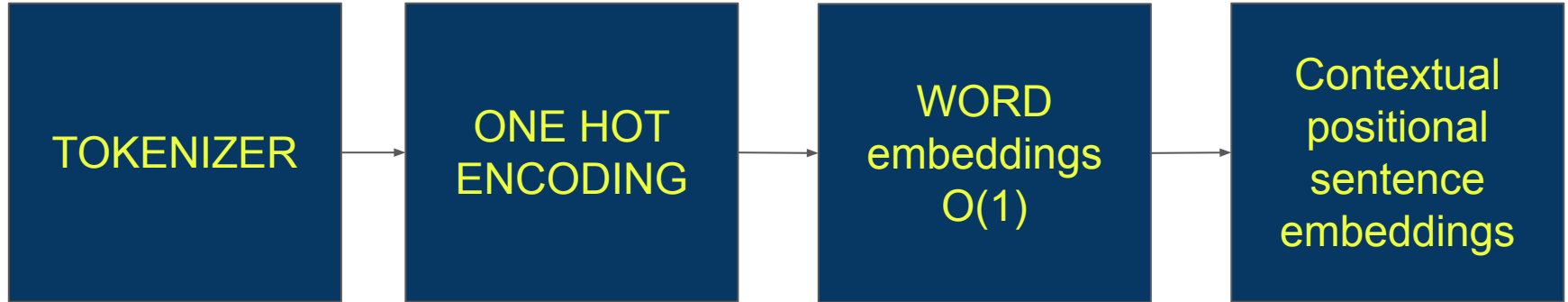
```
graph LR; A[TOKENIZER] --> B[ONE HOT ENCODING]; B --> C[WORD embeddings O(1)]
```

The diagram illustrates a three-step process for word representation. It begins with a 'TOKENIZER' block, which leads to a 'ONE HOT ENCODING' block, which in turn leads to a 'WORD embeddings O(1)' block. Each step is contained within a dark blue square, and the steps are connected by horizontal arrows pointing from left to right.

Contextual Word/Sentence embeddings

- Create a new linear transformation to transform word embeddings to produce
 - Contextual embeddings
 - Positional embeddings
 - Sentence embeddings

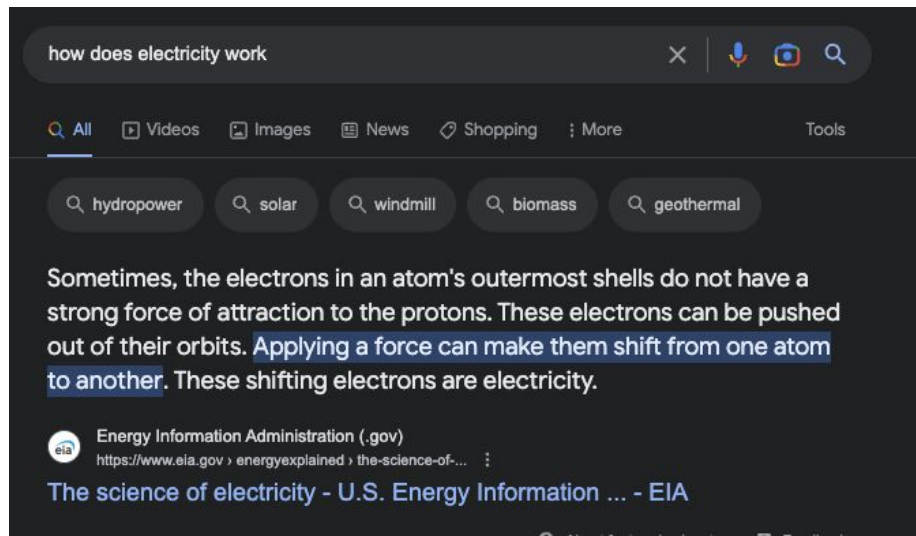
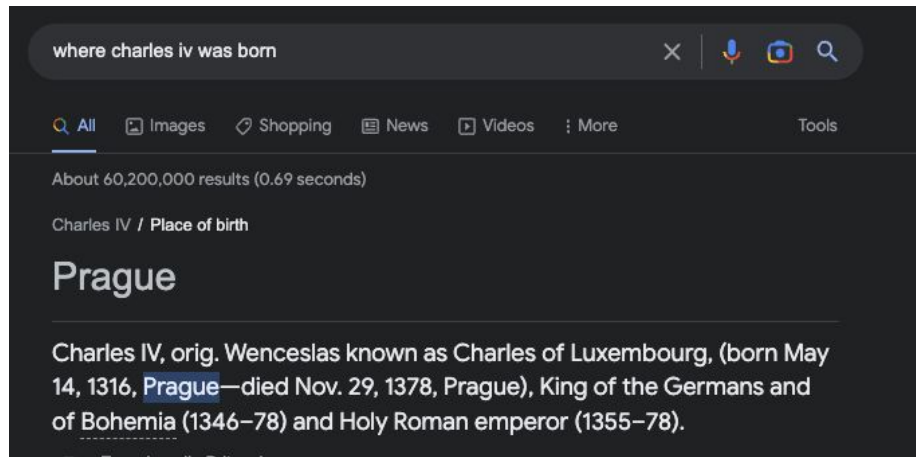




Fine-tune Embeddings

Fine tuning

- Where / how do we use embeddings?
- Fine tuning M input => 1 output
 - Google search
 - Question answering
 - Sentiment analysis
 - Toxicity classification
 - Automatic data annotation
- Compare similarity
 - Using vector similarity
 - “I like movies” ~ “cinemas are awesome”



Generative models

—— *with naive bayes* ——

M words on input \rightarrow N words on output

- Translators
- Summary text
- ChatGPT

GPT LLM (generative pretrained transformer)

Predict next word

Apply auto-regressive loop

Based on probabilities, space transformations may generate something new

Setup Temperature with naive bayes

Hello, how are you → my

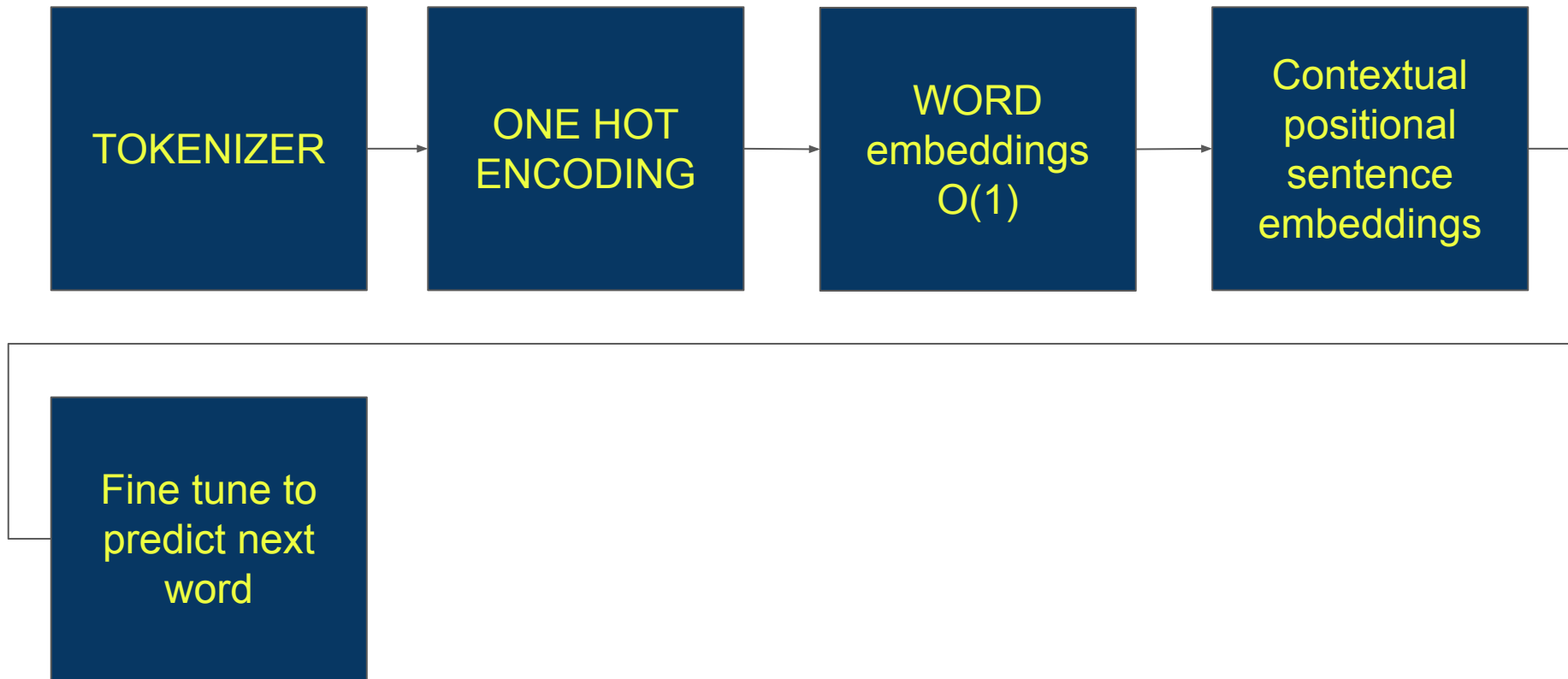
Hello, how are you my → friend

Hello, how are you my friend → ?

Bloomreach is my → favorite

Bloomreach is my favorite → company

Bloomreach is my favorite company → 🤖

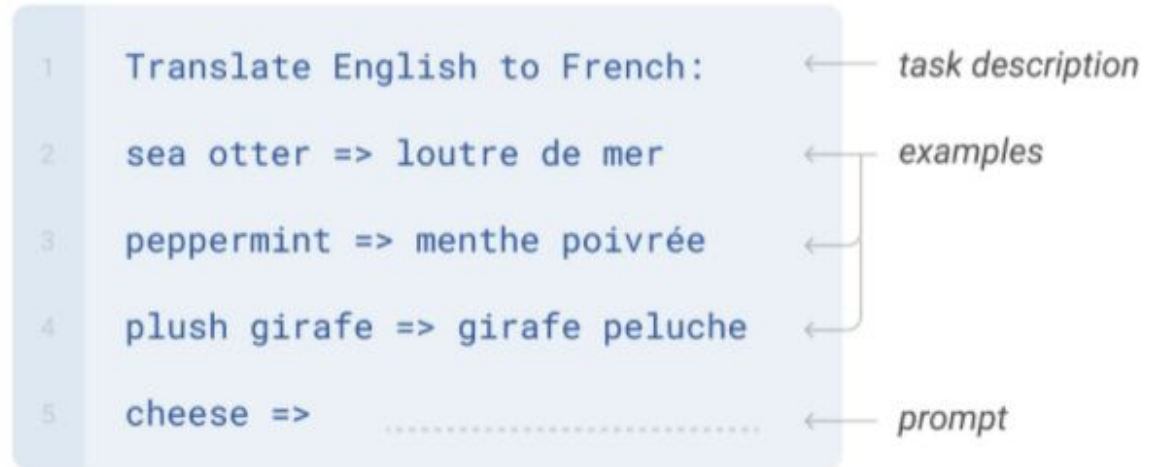


Fine-tune LLM in chatLGP

chatGPT *fine tuning*

Zero/one/few shot learning

Create a lot of chats and fine tune LLM to predict next word as it will be an conversation between 2 subjects



instruct GPT *fine tuning*

Prompt *Explain the moon landing to a 6 year old in a few sentences.*

Completion GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

InstructGPT

People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.

chat GPT *fine tuning*

Hello, how are you



→ my

Hello, how are you my



→ friend

Hello, how are you my friend



→ ?

Hello, how are you



→ I'm

Hello, how are you I'am



→ Just

Hello, how are you I'm Just



→ LLM

Human-curated AI safety *fine tuning*

Who is the best president?



→ I

Who is the best president? I

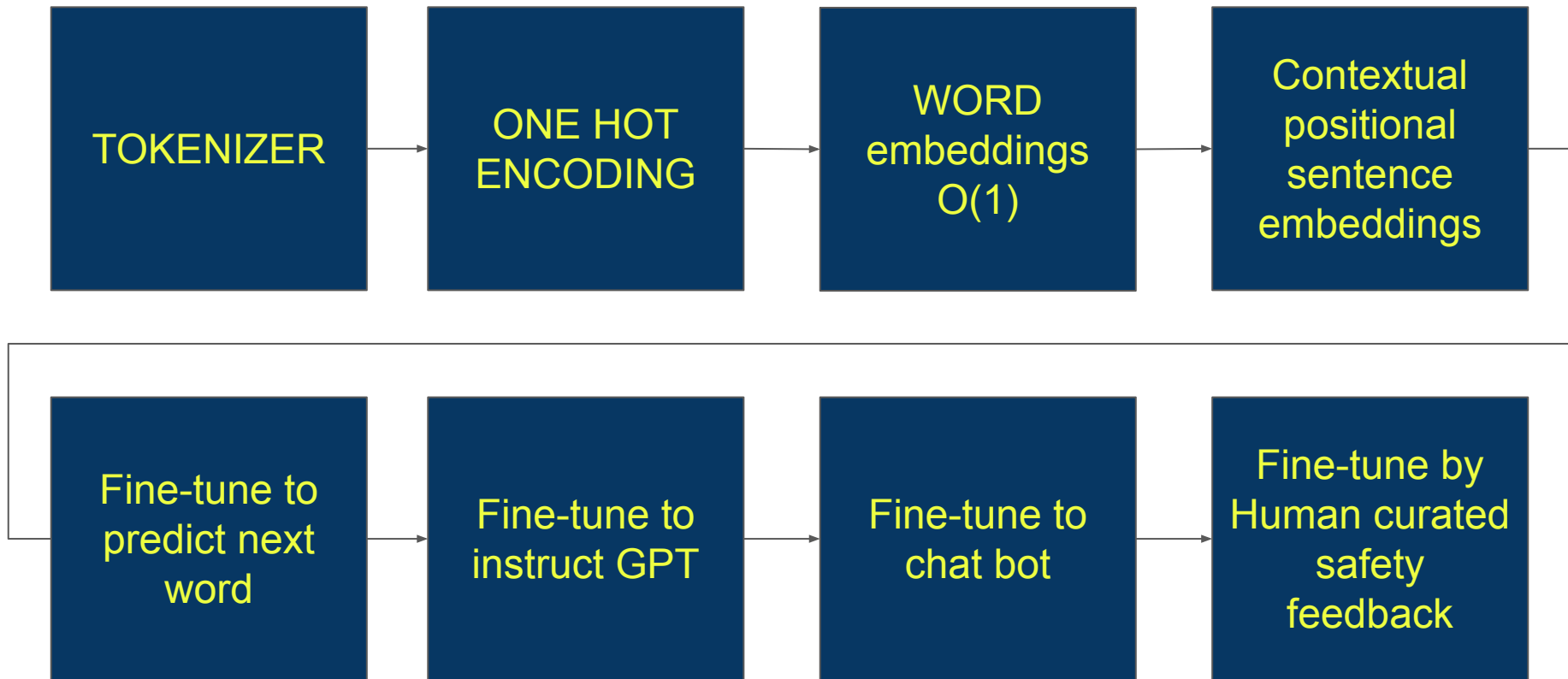


→ can't

Who is the best president? I can't



→ answer

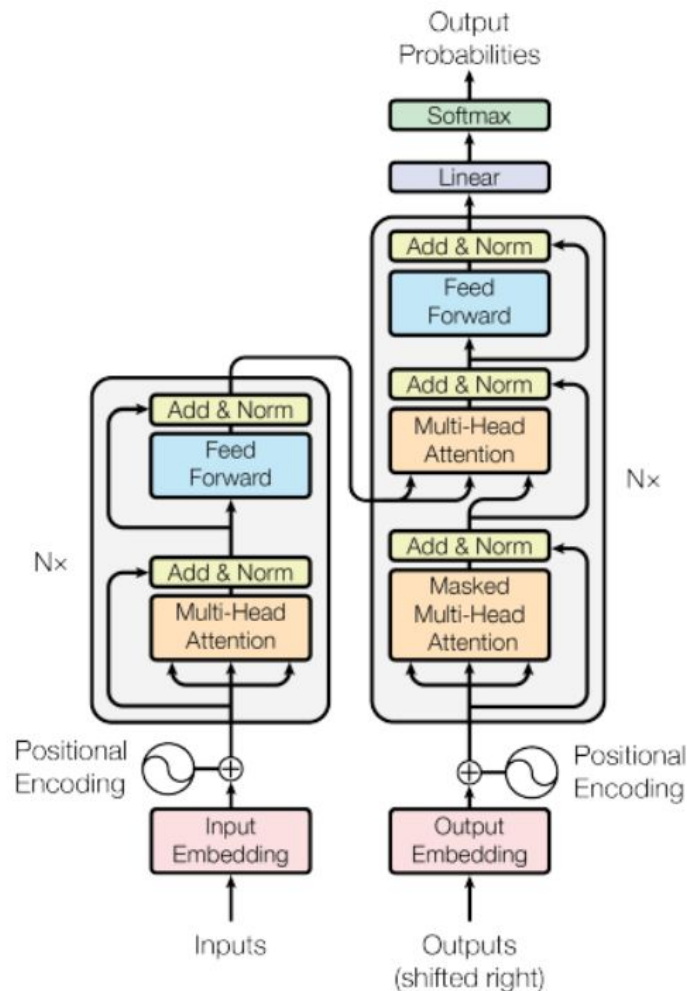


chatLGPT → chatGPT



chatLGPT → chatGPT

- Make it larger
- Make it faster
- Make it more distributable
- Make it larger
- Make it faster
- Add gigant corpus



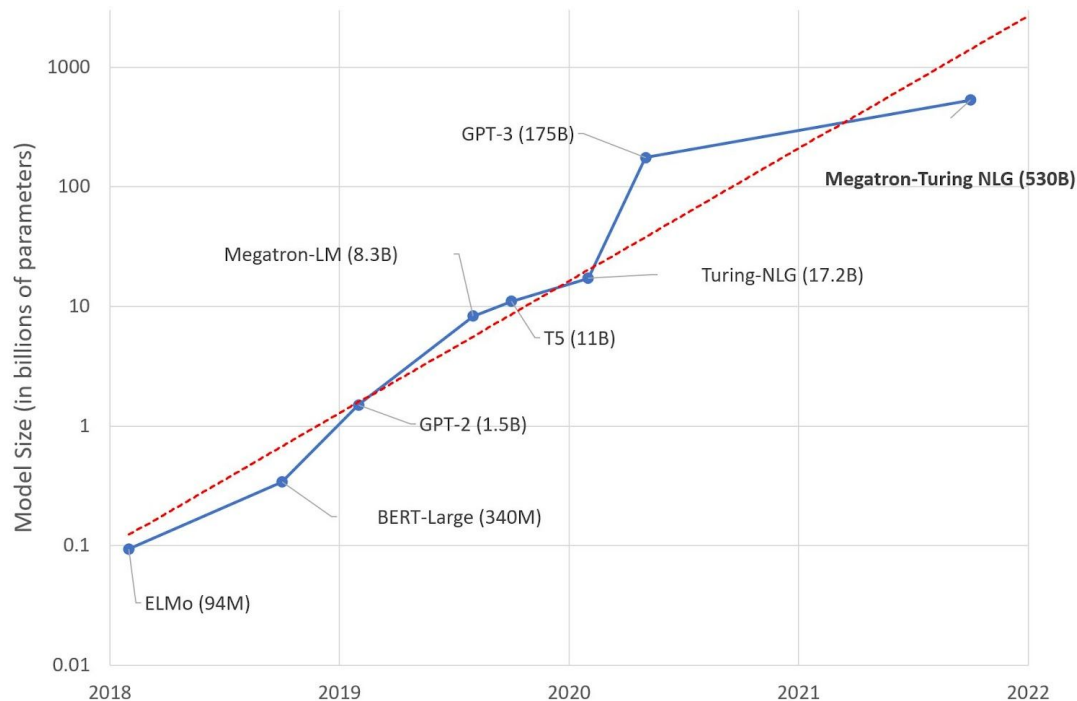
LLM comparison

LGPT 2023 (our)

- Corpus: 6 sentences
- Token count: 3
- parameters count: 14
- Window size: 2
- Embed vector dimension: 2

Open ai gpt3

- Corpus: multilingual internet + millions of books
- parameters count: 175B
- Window size: 2048
- Embed vector dimension: 12288



Mmm... who cares?

- LLM are “just” vector approximation machines
- LLM does not hold exact single source of truth
- LLM can't keep factual database
- LLM are not the tools which solve 100% of your use cases
- Knowledge & hallucinating of LLM is based on data which we use for learning

Current challenges in the field of NLP

- Short memory (people are trying to solve it with vector databases)
- $O(N^2)$ memory consumption (window size)
- Gigantic hardware requirements
- Needs to scrape the whole internet to do just dump model
- hallucinations
- It works the same as human think => word by word till the book exist
- Many AI innovations are just how to optimise learning to be able to train the model in finite time on current hardware
- How to use GPT?
 - It's not good to use it for fact-checking!!!
 - Its good for creative work
 - For example to summarize exact data which are factchecked like bing search do

Ezoteric future predictions

My humble future predictions

- FE UI will be dead => unified text interface to aggregate information across services (see bing chatbot)
- Code documentation like JSDoc is duplicity with the code itself
- Data types in programming are duplicates with runtime logic
- Documentation pages of programming libraries is duplicity with the code itself

- Natural languages gonna be less & less important because we have abstract vector
- LLM will be used as code generator for exact task where bugs are not allowed
- LLM will be used as interpreters for task which are not so exact

- Tutorials are dead, you only need to find a way how to put your raw data information into LLM
- Hard tradeoff between creating new nice API and generate old ugly but already included in the LLMs
- Monoliths vs microservices => it depends on window-size of LLM => current SOTA LLM has around 1K lines=> it could generate whole microservice

END - Disclaimer

- Examples were simplified into small neural network with ~10 params
- Real LLM NLP neural nets has which I did not cover here:
 - Batch normalization
 - Multi head (self-)Attention Layer
 - Positional encoding (i could add it into animation 🤔)
 - Residual connections (known from resNet)
 - Softmax (i could add it into animation 🤔)
 - dropout